

1 **hogwash: Three Methods for Genome-Wide Association Studies in Bacteria**  
2

3 **Authors**

4 Katie Saund<sup>1</sup> (0000-0002-6214-6713) and Evan S. Snitkin<sup>1,2</sup> (0000-0001-8409-278X)  
5

6 **Affiliations**

7 <sup>1</sup>Department of Microbiology and Immunology

8 <sup>2</sup>Department of Internal Medicine/Division of Infectious Diseases

9 University of Michigan, Ann Arbor, Michigan  
10

11 **Corresponding Author**

12 Evan S. Snitkin, esnitkin@med.umich.edu  
13

14 **Keywords**

15 GWAS, bacterial genomics, convergent evolution, software  
16

17 **ABSTRACT**

18 Bacterial genome-wide association studies (bGWAS) capture associations between genomic  
19 variation and phenotypic variation. Convergence based bGWAS methods identify genomic  
20 mutations that arise more often in the presence of phenotypic variation than is expected by  
21 chance. This work introduces hogwash, an open source R package that implements three  
22 algorithms for convergence based bGWAS. Hogwash additionally contains a novel grouping tool  
23 to perform gene- or pathway-analysis to improve power and increase convergence detection for  
24 related but weakly penetrant genotypes. To identify optimal use cases we applied hogwash to  
25 data simulated with a variety of phylogenetic signals and convergence distributions. These  
26 simulated data are publicly available and contain the relevant metadata regarding convergence  
27 and phylogenetic signal for each phenotype and genotype. Hogwash is available for download  
28 from GitHub.  
29

30 **DATA SUMMARY**

- 31 1. hogwash is available from GitHub under the MIT license  
32 (<https://github.com/katiesaund/hogwash>) and can be installed using the R command  
33 `devtools::install_github("katiesaund/hogwash")`
- 34 2. The simulated data used in this manuscript and the code to generate it are available  
35 from GitHub  
36 ([https://github.com/katiesaund/simulate\\_data\\_for\\_convergence\\_based\\_bGWAS](https://github.com/katiesaund/simulate_data_for_convergence_based_bGWAS))  
37

38 **IMPACT STATEMENT**

39 We introduce hogwash, an R package with three methods for bacterial genome-wide  
40 association studies. There are two methods for handling binary phenotypes, including an  
41 implementation of PhyC(1), as well as one method for handling continuous phenotypes. We  
42 formulate two novel indices quantifying the relationship between phenotype convergence and  
43 genotype convergence on a phylogenetic tree, one for binary phenotypes and one for  
44 continuous phenotypes. These indices shape an intuitive understanding for the ability of  
45 hogwash to detect significant intersections of phenotype convergence and genotype  
46 convergence and how to interpret hogwash outputs.  
47

48 **INTRODUCTION**

49 *Bacterial Genome-Wide Association Studies*

50 Bacterial genome-wide association studies (bGWAS) infer statistical associations between  
51 genotypes and phenotypes. Seminal bGWAS papers identified novel variants associated with

52 antibiotic resistance in *M. tuberculosis* and host specificity in *Campylobacter*(1,2). Since then,  
53 there have been numerous applications of bGWAS that have further highlighted the potential of  
54 this approach to identify genetic pathways underlying phenotypic variation and provide insights  
55 into the evolution of phenotypes of interest. Association studies can use various genetic data  
56 types including single nucleotide polymorphisms (SNPs), k-mers, copy number variants,  
57 accessory genes, insertions, and deletions. To improve the power and interpretability of bGWAS  
58 inclusion criteria or weighting can be applied to these variants based on predicted functional  
59 impact, membership in pathways of interest, or other user preferences(3,4). Differences  
60 between human and bacterial GWAS have been reviewed extensively by Power *et al.*(5). Of  
61 note, clonality and horizontal gene transfer complicate the application of human GWAS  
62 methodology to bacteria. However, certain bGWAS approaches can leverage unique features of  
63 bacterial evolution, including frequent phenotypic convergence and genotypic convergence, to  
64 identify phenotype-genotype correlations.  
65

#### 66 *bGWAS Software*

67 bGWAS methods can be classified into non-exclusive groups based on some critical features:  
68 A) methods for SNPs, accessory genes(6), or k-mers(7), B) methods using regression(7,8) or  
69 phylogenetic convergence(1,9), and C) methods designed for humans(10) or specifically for  
70 bacteria(7,9). Differences between regression based and convergence based bGWAS were  
71 expertly reviewed by Chen and Shapiro(11). Convergence based methods identify multiple  
72 independent events where a genomic mutation arises more often in the presence of the  
73 phenotype of interest. Convergence based methods can yield higher significance with a smaller  
74 sample size, but may fail to identify some statistical associations that traditional GWAS  
75 approaches would identify when the population is clonal(11). Additionally, convergence based  
76 methods are limited to smaller data sets because of their large memory requirements and  
77 computational time relative to traditional methods(12), but can surmount issues of clonality and  
78 take advantage of horizontal gene transfer.  
79

#### 80 *Objective*

81 This work describes hogwash, an R package that implements three different convergence  
82 based bGWAS approaches available on GitHub. Two approaches, PhyC and the Synchronous  
83 Test, handle binary phenotypes while the third approach, the Continuous Test, handles  
84 continuous phenotypes. PhyC is an algorithm introduced by Farhat *et al.*(1) that we implement.  
85 The Synchronous Test is a stringent variation of PhyC, requiring a tighter relationship between  
86 the genotype and phenotype. We describe the algorithms and evaluate them on a set of  
87 simulated data.  
88

#### 89 *Alternative Approach to Grouped Genotype Analysis*

90 Pathway analysis is a common post-GWAS approach that groups loci into meaningful groups,  
91 such as mapping SNPs to pathways(13,14). Analyzing aggregated loci can improve both the  
92 interpretability of GWAS results and improve power to detect associations(13,14). However,  
93 post-GWAS pathway analysis may not surmount the high stringency of convergence based  
94 bGWAS approaches. Hogwash implements a novel grouping tool prior to performing the  
95 bGWAS that may avoid this potential loss of information and improve convergence detection for  
96 related but weakly penetrant genotypes.  
97

#### 98 *Data Simulation*

99 We evaluate hogwash results on simulated data generated to capture aspects of bacterial  
100 evolution pertinent to these bGWAS approaches. We simulated data with a range of  
101 phylogenetic signals and convergence distributions to highlight the critical impact of these  
102 features on bGWAS results. The simulated data are publicly available and could be used to

103 compare the impact of convergence patterns within phenotypes, genotypes, and their  
104 intersection when benchmarking various convergence based bGWAS methods.  
105

## 106 PACKAGE DESCRIPTION

107 We developed hogwash to allow users to perform three bGWAS methods, including an open  
108 source implementation of the previously described PhyC algorithm(1), and aggregate genotypes  
109 by user-defined groups. The hogwash function minimally requires a phenotype, a phylogenetic  
110 tree, and a binary genotype matrix. An optional argument may be supplied to facilitate grouping  
111 genotypes. The genotype matrix and tree can be prepared from a multiVCF file by the variant  
112 preprocessing tool prewas(15). Hogwash assumes that the genotype is encoded such that 0  
113 refers to wild type and 1 refers to a mutation and that binary phenotypes are encoded such that  
114 0 refers to absence and 1 refers to presence.  
115

116 In brief, the hogwash workflow (Figure 1A) begins with the user supplying a phenotype, a set of  
117 genotypes, and a tree. Hogwash performs ancestral state reconstruction for the phenotype and  
118 genotypes. If the user supplies a key to group together genotypes, hogwash groups them after  
119 the genotype ancestral state reconstructions. The convergence of each phenotype and each  
120 genotype are recorded as the edges where they intersect on the tree (Figure 1B); the definition  
121 of convergence and intersection is unique for each of the three association tests. Then the  
122 genotype is permuted, and its intersection with the phenotype is recorded as a null distribution.  
123 Significance is calculated with correction for multiple testing.  
124

## 125 PhyC

126 PhyC is a convergence based bGWAS method introduced by Farhat *et al.*(1) that identified  
127 novel antibiotic resistance-conferring mutations in *M. tuberculosis*. To our knowledge, the  
128 original PhyC code is not publicly available, but the algorithm is well described in the original  
129 paper. The algorithm addresses the following question: Does the genotype transition from wild  
130 type, 0, to mutant, 1, occur more often than expected by chance on tree edges where the  
131 phenotype is present, 1, than where the phenotype is absent, 0? By requiring the overlap of the  
132 phenotype with the genotype transition, instead of genotype presence, associations are not  
133 inflated by clonal sampling and thus this approach controls for population structure. We  
134 implement the PhyC algorithm as described in Farhat *et al.*(1).  
135

136 For each test we formulate the following terms:  $\beta_{genotype}$  and  $\beta_{phenotype}$ . In PhyC (Figure 2)  
137  $\beta_{genotype}$  is a non-negative integer that records the number of tree edges where the genotype  
138 arises (mutation appears). This is encoded as a tree edge where the parent node is evaluated  
139 as 0 by ancestral state reconstruction and child node is evaluated as 1. These edges are called  
140 genotype transitions. In PhyC  $\beta_{phenotype}$  is a non-negative integer that records the number of  
141 tree edges where the phenotype is present. This is encoded as a tree edge where the child  
142 node is evaluated as 1 by ancestral state reconstruction. The number of edges on the tree  
143 where both a genotype arises and the phenotype is present is calculated as  $\beta_{genotype} \cap$   
144  $\beta_{phenotype}$ .  
145

146 For the permutation the genotype mutations ( $\beta_{genotype}$ ) are randomized on the tree. The  
147 number of edges where the permuted genotype mutation intersects with phenotype presence  
148 edges is recorded for each permutation; these permuted  $\beta_{genotype} \cap \beta_{phenotype}$  values create a  
149 null distribution. An empirical *P*-value is calculated based on the observed  $\beta_{genotype} \cap \beta_{phenotype}$   
150 as compared to the null distribution.  
151

152 Our PhyC implementation has several important differences from the original paper. First,  
153 multiple test correction in hogwash is performed with False Discovery Rate instead of the more  
154 stringent Bonferroni correction. Second, hogwash reduces the multiple testing burden by testing  
155 only those genotype-phenotype pairs for which convergence is detectable; genotypes with  
156  $\beta_{genotype} < 2$  are excluded and genotype-phenotype pairs with  $\beta_{genotype} \cap \beta_{phenotype} < 2$  are  
157 assigned a  $P$ -value of 1. Third, ancestral state reconstruction for genotypes and phenotypes  
158 was performed using only maximum likelihood. Finally, users only supply one phylogenetic tree  
159 instead of three.  
160

### 161 **Synchronous Test**

162 This test (Figure 2) is an extension of PhyC but requires more stringent association between the  
163 genotype and phenotype. The Synchronous Test addresses the question: Do genotype  
164 transitions occur more often than expected by chance on phenotype transition edges than on  
165 phenotype non-transition edges? Both PhyC and the Synchronous Test are only appropriate for  
166 binary phenotypes.  
167

168 The Synchronous Test  $\beta_{genotype}$  is a non-negative integer that records the number of tree edges  
169 where the genotype changes (mutation appears or disappears). This is encoded as a tree edge  
170 where the parent node value as inferred from ancestral state reconstruction is different than the  
171 child node. These edges are called genotype transitions. The Synchronous Test  $\beta_{phenotype}$  is a  
172 non-negative integer that records the number of tree edges where the phenotype changes. This  
173 is encoded as a tree edge where the phenotype parent node is different than the child node as  
174 inferred from ancestral state reconstruction. The number of edges on the tree where both a  
175 genotype transitions and the phenotype transitions is calculated as  $\beta_{genotype} \cap \beta_{phenotype}$ . As in  
176 PhyC, the genotypes with  $\beta_{genotype} < 2$  are removed, genotype-phenotype pairs with  
177  $\beta_{genotype} \cap \beta_{phenotype} < 2$  are assigned a  $P$ -value of 1, and the remaining genotypes are  
178 permuted and a null distribution of the  $\beta_{genotype} \cap \beta_{phenotype}$  is calculated to determine the  
179 significance of each genotype.  
180

181 This test is similar to the Simultaneous Score in treeWAS(9). The Simultaneous Score is  
182 derived from the number of edges on the tree where the genotype and phenotype transition in  
183 the same direction (both have an inferred parent node of 0 and an inferred child node of 1 or  
184 parent node of 0 and child node of 1). In contrast, the Synchronous Test in hogwash allows for  
185 the phenotype and genotype transition directions to mismatch, thus allowing for the detection of  
186 genotypes with inconsistent effect directions.  
187

### 188 **Continuous Test**

189 The Continuous Test (Figure 2) is an application of a convergence based GWAS method to  
190 continuous phenotypes. The Continuous Test addresses the question: Does the phenotype  
191 change more than expected by chance on genotype transition edges than on genotype non-  
192 transition edges?  
193

194 The Continuous Test  $\beta_{genotype}$  is a non-negative integer that records the number of tree edges  
195 where the genotype changes (mutation appears or disappears). This is encoded as a tree edge  
196 where the parent node value as inferred from ancestral state reconstruction is different than the  
197 child node (Figure 1B). These edges are called genotype transitions.  
198

### **Formula 1.**

$$\Delta_{edge} = |phenotype_{parent\ node} - phenotype_{child\ node}|$$

199

200 A  $\Delta_{edge}$  value is calculated for each tree edge and is scaled from 0 to 1. The Continuous Test  
201  $\beta_{phenotype}$  is the sum of all  $\Delta_{edge}$  values that occur on high confidence edges.  $\beta_{genotype} \cap$   
202  $\beta_{phenotype}$  is the multiplicative sum of the  $\Delta_{edge}$  and  $\beta_{genotype}$ . As above, the genotypes with  
203  $\beta_{genotype} < 2$  are removed; the remaining genotypes are permuted and a null distribution of  
204 the  $\beta_{genotype} \cap \beta_{phenotype}$  is calculated to determine the significance of each genotype.  
205

## 206 **User inputs**

207 The user must provide a phylogenetic tree, genotype matrix, and a phenotype. The user may  
208 optionally provide a key that maps individual genomic loci into groups in order to use hogwash's  
209 grouping feature. For a detailed description of the user inputs please see the Supplementary  
210 Package Description.

211

## 212 **Hogwash outputs**

213 The package produces two files per test: data (.rda) and plots (.pdf). The data file contains  
214 many pieces of information, including  $P$ -values for each tested genotype. The plots are  
215 described below in the results section.

216

## 217 **Grouping feature**

218 To identify an association between a genomic variant and a phenotype, hogwash requires that a  
219 variant occur in multiple different lineages. Hogwash may classify some causal variants as  
220 independent of a phenotype if they are weakly penetrant, an issue common to convergence  
221 based methods. To surmount this issue, related genomic variants may be aggregated to capture  
222 larger trends at the grouped level. For example, a user may apply this method to group only  
223 nonsynonymous SNPs by gene to use hogwash to detect associations between the mutated  
224 gene and the phenotype. Grouping related variants can improve power through a reduction in  
225 the multiple testing correction penalty. However, the power benefits are dependent on grouping  
226 variants with similar effect directions.

227

228 Hogwash implements the grouping features by first performing ancestral state reconstruction for  
229 each individual locus (Figure 3). If the user supplies a key that maps individual loci to groups,  
230 then the edges contributing to  $\beta_{genotype}$  for individual loci are joined together into the indicated  
231 group. Grouped loci with  $\beta_{genotype} < 2$  are excluded from analysis. After this point hogwash  
232 runs as previously described for non-grouped genotypes.

233

234 Users may supply hogwash with data that was previously grouped (for example, using the group  
235 SNPs by gene functionality in prewas(15)) but this approach may mask some genotype  
236 transitions. In this case, the user does not need to provide a key and the hogwash grouping step  
237 is skipped.

238

## 239 **METHODS**

### 240 **Data simulation**

#### 241 *Trees*

242 We simulated four random coalescent phylogenetic trees with 100 tips each.

243

#### 244 *Phenotypes*

245 For each tree we simulated phenotypes that model either Brownian motion or white noise. A  
246 phenotype modeled well with white noise may suggest a role for horizontal gene transfer, gene  
247 loss, or convergent evolution(16). A white noise phenotype may be better suited to the hogwash  
248 algorithms than a phenotype modeled by Brownian motion given the requirement for  
249 phylogenetic convergence. A continuous phenotype that is modeled well by Brownian motion

250 has a phylogenetic signal,  $\lambda$ , near 1 while a white noise phenotype has a phylogenetic signal  
251 near 0(17). In contrast, a binary phenotype that is modeled well by Brownian motion has a  
252 phylogenetic signal,  $D$  statistic, near 0 while a white noise phenotype has a phylogenetic signal  
253 near 1(18). For each tree we simulated sixteen phenotypes: eight phenotypes with a  
254 phylogenetic signal fitting a Brownian motion model (four binary and four continuous), and eight  
255 phenotypes with a phylogenetic signal fitting a white noise model (four binary and four  
256 continuous). For phenotypes modeling Brownian motion, binary phenotypes were restricted to  
257  $-0.05 < D < 0.05$  and continuous phenotypes to  $0.95 < \lambda < 1.05$ . For phenotypes modeling  
258 white noise, binary phenotypes were restricted to  $0.95 < D < 1.05$  and continuous phenotypes  
259 to  $-0.05 < \lambda < 0.05$ .

260

### 261 **Genotypes**

262 For each simulated tree a set of unique binary genotypes were generated. We generated  
263 genotypes that span a range of phylogenetic signals, degree of similarity to the phenotype, and  
264 prevalence.

#### 265 Genotypes to be used in discrete hogwash tests

266 First, 25,000 binary genotypes were generated using ape::rTraitDisc; these genotypes have a  
267 range of phylogenetic signals(19). Second, these genotypes were duplicated and randomized  
268 with the following approach: one quarter had 10% of tips changed, one quarter had 25% of tips  
269 changed, one quarter had 40% of tips changed, and one quarter were entirely redistributed.  
270 Third, we removed any simulated genotypes present in 0, 1,  $N - 1$ , or  $N$  samples. Fourth, we  
271 subset the genotypes to keep only unique presence/absence patterns. Fifth, we subset  
272 genotypes to only those within a range of  $-1.5 < D < 1.5$ . These filtering steps result in a  
273 reduction in the data set size (range 2214-2334).

#### 274 Genotypes to be used in the Continuous Test

275 In addition to the five steps above we added two more data generation steps. First, we made all  
276 possible genotypes based on the rank of the continuous phenotype. Second, we made  
277 genotypes based on which edges of the tree had high  $\Delta_{edge}$ . The filtering steps reduced the  
278 data set size (range 1234-1310).

279

### 280 **Hogwash on simulated data**

281 We ran each hogwash test for each of the tree-phenotype-genotype sets. In addition to  
282 generating  $P$ -values for each tested genotype, hogwash also reports convergence information.  
283 We ran hogwash with the following settings: permutations = 50,000; false discovery rate =  
284 0.0005 (discrete), 0.05 (continuous); bootstrap value = 0.70; no genotype grouping key was  
285 provided.

286

#### 287 Calculation of $\beta_{genotype}$ , $\beta_{phenotype}$ and $\varepsilon$

288 Using the ancestral state reconstruction data, hogwash identified convergence within each  
289 genotype ( $\beta_{genotype}$ ), phenotype ( $\beta_{phenotype}$ ), and their weighted intersection ( $\varepsilon$ ).

#### 290 **Formula 2.**

$$\varepsilon_{binary} = \frac{2 \times (\beta_{genotype} \cap \beta_{phenotype})}{\beta_{genotype} + \beta_{phenotype}}$$

291

292 Where  $0 \leq \beta_{genotype} \leq Number\ tree\ tips$  and  $0 \leq \beta_{phenotype} \leq Number\ tree\ tips$ ; both  
293  $\beta_{genotype}$  and  $\beta_{phenotype}$  are integers. Therefore,  $0 \leq \varepsilon \leq 1$ .

294

#### 295 **Formula 3.**

$$\varepsilon_{continuous} = \frac{\beta_{genotype} \cap \beta_{phenotype}}{\beta_{genotype} + \beta_{phenotype} - \beta_{genotype} \cap \beta_{phenotype}}$$

296

297 Where  $0 \leq \beta_{genotype} \leq Number\ tree\ tips$  and  $0 \leq \beta_{phenotype} < Number\ tree\ tips$ .  $\beta_{genotype} \cap \beta_{phenotype}$  is the multiplicative sum of the  $\Delta_{edge}$  and  $\beta_{genotype}$ . The denominator is  $\beta_{genotype} \cup \beta_{phenotype}$ . Therefore,  $0 \leq \varepsilon \leq 1$ .

300

### 301 Data analysis

302 Statistical analyses were conducted in R v3.6.2(20). The R packages used can be found in the  
303 simulate\_data.yaml file on GitHub(19,21–25) and can be installed using miniconda(26).

304

## 305 RESULTS

### 306 Hogwash output for simulated data

307 Hogwash outputs two sets of results: a data file and a PDF file with plots. Each run of PhyC  
308 produces at least three plots: the phenotype reconstruction (Figure 4A), a Manhattan plot  
309 (Figure 4D), and a heatmap of all tested genotypes (Figure 4E). The phenotype reconstruction,  
310 also referred to as  $\beta_{phenotype}$ , is highlighted on the tree (Figure 4A). The Manhattan plot shows  
311 the distribution of  $P$ -values from the hogwash run (Figure 4D). Lastly, the heatmap shows the  
312 genotype reconstruction ( $\beta_{genotype}$ ) and phenotype reconstruction ( $\beta_{phenotype}$ ) for each tree  
313 edge (rows) and genotype (columns) (Figure 4E). The genotypes are clustered by the  $\beta_{genotype}$   
314 presence/absence pattern. Two additional plots are produced for each genotype that is  
315 significantly associated with the phenotype: a phylogenetic tree showing the genotype transition  
316 edges (Figure 4B) and the null distribution of  $\beta_{genotype} \cap \beta_{phenotype}$  (Figure 4C).

317

318 The Synchronous Test and Continuous Test output plots that reflect their test-specific  
319  $\beta_{genotype}$  and  $\beta_{phenotype}$  definitions (Figure S1, S2). Running hogwash on 100 samples required  
320 <3 hours and <2 GB of memory for binary data and <5 hours and <2 GB of memory for  
321 continuous data (Figure S3).

322

### 323 Hogwash evaluation on simulated data

324 To help users identify optimal use cases and also interpret hogwash results we describe the  
325 behavior of hogwash on simulated data. We note that this assessment is not meant to convey  
326 performance in the sense of calculating sensitivity and specificity, but rather evaluate whether  
327 hogwash can robustly detect the association between phenotypic and genotypic convergence.  
328 To guide our assessment, we compared the relationship between the  $P$ -value and  $\varepsilon$  values  
329 produced by hogwash on sets of simulated data constructed using different evolutionary models  
330 (Figure 5).  $\varepsilon$  is a quantification of the relationship between phenotype convergence and  
331 genotype convergence; we define  $\varepsilon$  for the discrete and continuous tests in Formulae 2 and 3,  
332 respectively. Low  $\varepsilon$  values indicate little to no intersection of phenotype convergence and  
333 genotype convergence, while higher  $\varepsilon$  values indicate their increased intersection.

334

335 For discrete phenotypes, we observe an overall strong positive association between  $-\log(P$ -  
336 value) and  $\varepsilon$ , demonstrating that as the intersection of phenotype convergence and genotype  
337 convergence increase hogwash predicts that it is less likely that they intersect due to chance  
338 (Table 1). In other words, below a certain  $\varepsilon_{binary}$  threshold, hogwash attributes the association  
339 between the genotype convergence and phenotype convergence to chance; from Figure 5 the  
340 user can get a sense for the range of this  $\varepsilon_{binary}$  threshold under different evolutionary regimes.

341

342 For the simulated continuous data an  $\varepsilon_{continuous}$  threshold that separates meaningful genotype-  
343 phenotype associations from associations by chance is less apparent. Higher  $\varepsilon$ , low significance  
344 values demonstrate that some overlap of  $\beta_{genotype}$  and  $\beta_{phenotype}$  is likely by chance given the  
345 data. Low  $\varepsilon$ , high significance values demonstrate that some values with even small amounts of  
346  $\beta_{genotype}$  and  $\beta_{phenotype}$  overlap are unlikely, however that does not necessarily suggest that  
347 these hits are the best candidates for *in vitro* follow up. We suspect that these associations are  
348 largely driven by poor exploration of the sampling space, despite running many permutations,  
349 because of the edge-length based sampling probability of the permutation method. Therefore, it  
350 is essential  $P$ -values be interpreted within the context of  $\varepsilon$ . Notably, the Continuous Test was  
351 only able to detect significant genotype-phenotype associations for phenotypes modeled by  
352 white noise, suggesting this method is particularly sensitive to the phenotype's evolutionary  
353 model.

354

355 We observe for both the discrete and continuous tests that  $\varepsilon$  is more tightly correlated with -  
356  $\log(P\text{-value})$  for phenotypes characterized by white noise than by Brownian motion (Table 1),  
357 indicating that hogwash performs better under a white noise model. Thus, to help the user in  
358 assessing the appropriateness of hogwash and in interpreting their results we allow users to  
359 check if their phenotype is better modeled by white noise than Brownian motion by using the  
360 report\_phylogenetic\_signal function.

361

362

## 363 DISCUSSION

364 We have developed hogwash, an open-source R package that implements three different  
365 approaches to bGWAS and includes the previously described PhyC algorithm(1). Hogwash also  
366 introduces a novel grouping feature to aggregate related genomic variants to increase detection  
367 of convergence for weakly penetrant genotypes. Hogwash is best used for datasets comprising  
368 binary and/or continuous phenotypes, phenotypes fitting white noise models, situations where  
369 convergence may occur at the level of genes or pathways and with datasets whose size can be  
370 accommodated given the time and memory constraints of convergence methods.

371

372 The results of running hogwash on simulated data suggest that after a certain  $\varepsilon$  threshold, it  
373 unlikely that the intersection between phenotype convergence and genotype convergence  
374 occurs by chance, particularly for white noise phenotypes. Given the variability in results within  
375 each method, as shown in Figure 5, users may want to contextualize the statistical significance  
376 of the tested genetic loci with the amount of convergence possible for any one particular data  
377 set; to facilitate this the hogwash output includes both  $P$ -values and  $\varepsilon$ .

378

379 The simulated data set presented here is published to serve as a resource or template for future  
380 work focused on benchmarking convergence based bGWAS software as such a dataset has not  
381 yet, as far as we are aware, been made available(27). The simulated data set is available on  
382 GitHub and includes convergence information for each phenotype, genotype, and their  
383 intersection.

384

## 385 AUTHOR CONTRIBUTIONS

386 KS and ESS conceptualized the project and edited the manuscript. KS designed and  
387 implemented the software, performed the analysis, prepared the original draft, and visualized  
388 the data. ESS supervised the project.

389

## 390 CONFLICTS OF INTEREST

391 The authors declare that there are no conflicts of interest.

392

## 393 FUNDING

394 KS was supported by the National Institutes of Health (T32GM007544). ESS and KS were  
395 supported by the National Institutes of Health (1U01AI124255).

396

## 397 ACKNOWLEDGEMENTS

398 We thank Brad Saund for his help formalizing the continuous algorithm  $\varepsilon$  definition.

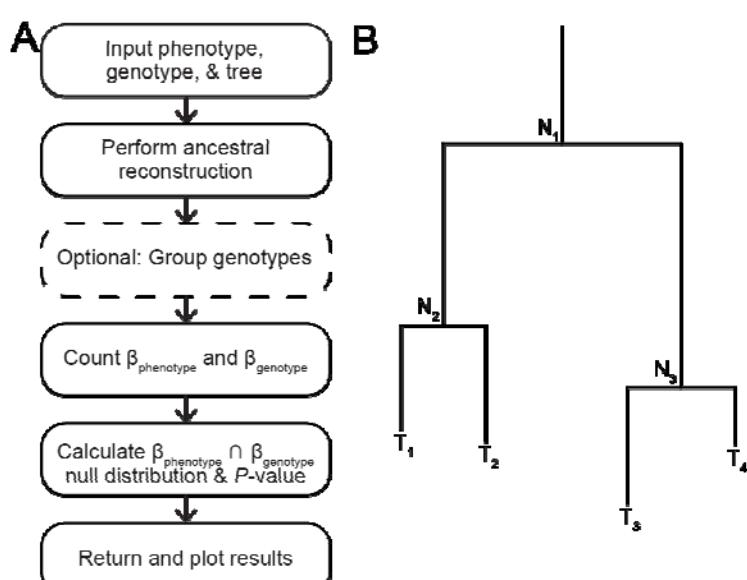
399

## 400 REFERENCES

- 401 1. Farhat MR, Shapiro BJ, Kieser KJ, Sultana R, Jacobson KR, Victor TC, et al. Genomic  
402 analysis identifies targets of convergent positive selection in drug-resistant  
403 Mycobacterium tuberculosis. *Nat Genet.* 2013;45(10):1183–9.
- 404 2. Sheppard SK, Didelot X, Meric G, Torralbo A, Jolley KA, Kelly DJ, et al. Genome-wide  
405 association study identifies vitamin B 5 biosynthesis as a host specificity factor in  
406 *Campylobacter*. *PNAS.* 2013;110(29):11923–7.
- 407 3. Sveinbjornsson G, Albrechtsen A, Zink F, Gudjonsson SA, Oddson A, Másson G, et al.  
408 Weighting sequence variants based on their annotation increases power of whole-  
409 genome association studies. *Nat Genet.* 2016 Mar 1;48(3):314–7.
- 410 4. Hendricks AE, Bochukova EG, Marenne G, Keogh JM, Atanassova N, Bounds R, et al.  
411 Rare Variant Analysis of Human and Rodent Obesity Genes in Individuals with Severe  
412 Childhood Obesity. *Sci Rep.* 2017 Dec 1;7(1):1–14.
- 413 5. Power RA, Parkhill J, de Oliveira T. Microbial genome-wide association studies: lessons  
414 from human GWAS. *Nat Rev Genet.* 2016;
- 415 6. Brynildsrød O, Bohlin J, Scheffer L, Eldholm V. Rapid scoring of genes in microbial pan-  
416 genome-wide association studies with Scoary. *Genome Biol.* 2016;17.
- 417 7. Lees JA, Galardini M, Bentley SD, Weiser JN, Corander J. pyseer: a comprehensive tool  
418 for microbial pangenome-wide association studies. *Bioinformatics [Internet].* 2018 [cited  
419 2018 Dec 19]; Available from: <http://pyseer.readthedocs.io>.
- 420 8. Earle SG, Wu C-H, Charlesworth J, Stoesser N, Gordon NC, Walker TM, et al. Identifying  
421 lineage effects when controlling for population structure improves power in bacterial  
422 association studies. *Nat Microbiol.* 2016;1.
- 423 9. Collins C, Didelot X. A phylogenetic method to perform genome-wide association studies  
424 in microbes that accounts for population structure and recombination. *PLoS Comput Biol.*  
425 2018;
- 426 10. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al. PLINK: A  
427 Tool Set for Whole-Genome Association and Population-Based Linkage Analyses. *Am J  
428 Hum Genet [Internet].* 2007 [cited 2017 Mar 22];81(3):559–75. Available from:  
429 [www.ajhg.org](http://www.ajhg.org)
- 430 11. Chen PE, Shapiro BJ. The advent of genome-wide association studies for bacteria. Vol.  
431 25, *Current Opinion in Microbiology.* 2015. p. 17–24.
- 432 12. Corander J, Croucher NJ, Harris SR, Lees JA, Tonkin-Hill G. Bacterial Population  
433 Genomics. In: *Handbook of Statistical Genomics.* Wiley; 2019. p. 997–1020.
- 434 13. Mooney MA, Wilmot B. Gene set analysis: A step-by-step guide. *Am J Med Genet Part B  
435 Neuropsychiatr Genet.* 2015 Oct 1;168(7):517–27.
- 436 14. White MJ, Yaspan BL, Veatch OJ, Goddard P, Rissee-Adams OS, Contreras MG.  
437 Strategies for Pathway Analysis Using GWAS and WGS Data. *Curr Protoc Hum Genet.*  
438 2019 Jan 1;100(1):e79.
- 439 15. Saund K, Lapp Z, Thiede SN, Pirani A, Snitkin ES. prewas: Data pre-processing for more  
440 informative bacterial GWAS. *bioRxiv.* 2019 Dec 20;2019.12.20.873158.
- 441 16. van Assche A, Alvarez-Perez S, de Breij A, de Brabanter J, Willems KA, Dijkshoorn L, et  
442 al. Phylogenetic signal in phenotypic traits related to carbon source assimilation and

- 443 chemical sensitivity in *Acinetobacter* species. *Appl Microbiol Biotechnol*. 2016;101:367–  
444 79.  
445 17. Pagel M. Inferring the historical patterns of biological evolution. *Nature*.  
446 1999;401(6756):877–84.  
447 18. Fritz SA, Purvis A. Selectivity in mammalian extinction risk and threat types: A new  
448 measure of phylogenetic signal strength in binary traits. *Conserv Biol*. 2010;24(4):1042–  
449 51.  
450 19. Paradis E, Schliep K. Phylogenetics ape 5.0: an environment for modern phylogenetics  
451 and evolutionary analyses in R.  
452 20. R Core Team. R: A language and environment for statistical computing. R Foundation for  
453 Statistical Computing, Vienna, Austria.; 2018.  
454 21. Orme D. The caper package: comparative analysis of phylogenetics and evolution in R.  
455 R Packag version 05, 2. 2013;1–36.  
456 22. Revell LJ. phytools: An R package for phylogenetic comparative biology (and other  
457 things). *Methods Ecol Evol*. 2012;3(2):217–23.  
458 23. Wickham H. tidyverse: Easily Install and Load the “Tidyverse.” 2017.  
459 24. Wickham H, Seidel D. scales: Scale Functions for Visualization. 2019.  
460 25. Auguie B. gridExtra: Miscellaneous Functions for “Grid” Graphics. [Internet]. 2017.  
461 Available from: <https://cran.r-project.org/package=gridExtra>  
462 26. Anaconda [Internet]. [cited 2020 Feb 21]. Available from: <https://www.anaconda.com/>  
463 27. Saber MM, Shapiro BJ. Benchmarking bacterial genome-wide association study methods  
464 using simulated genomes and phenotypes. *Microb genomics*. 2020;6(3).

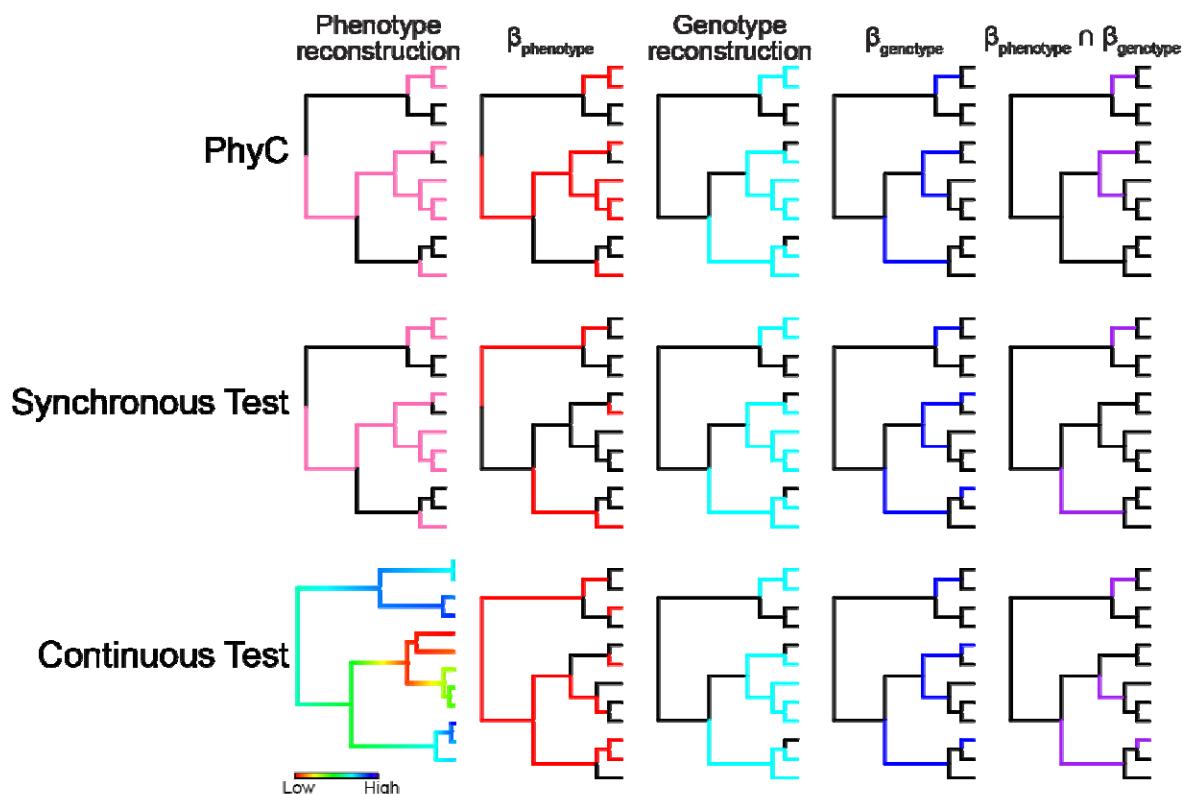
465  
466  
467 **FIGURES**  
468



469  
470  
471  
472  
473  
474  
475

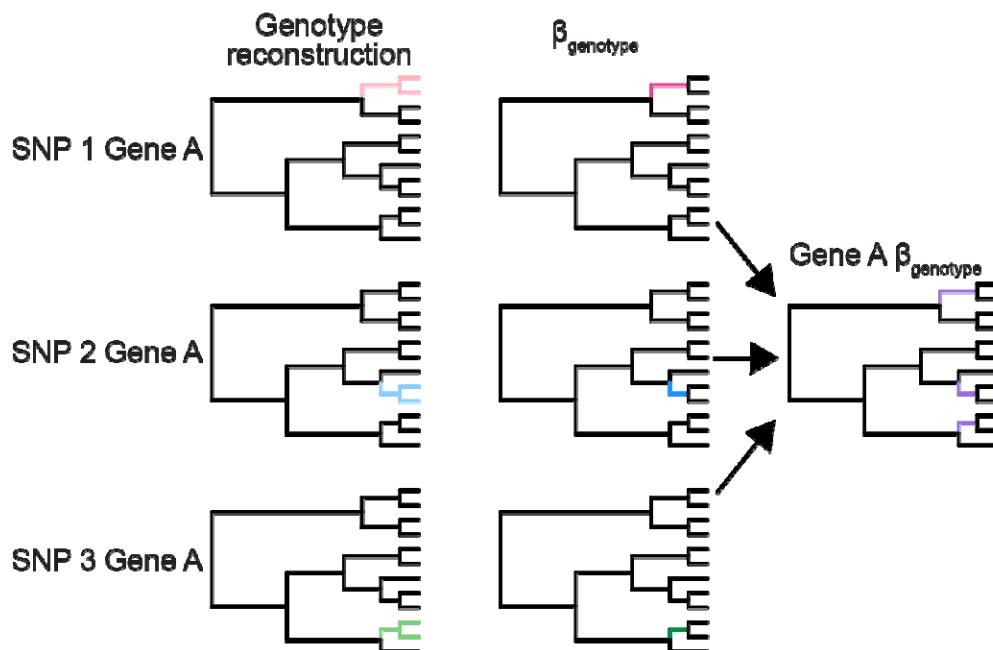
**Figure 1. Hogwash workflow and tree nomenclature.**

A) Software workflow. B) In this example phylogenetic tree N<sub>1</sub> is the root. Tree nodes are labeled N<sub>1</sub> – N<sub>3</sub>. Tree tips are labeled T<sub>1</sub> – T<sub>4</sub>. N<sub>1</sub> is a parent node to N<sub>2</sub> and N<sub>3</sub>. N<sub>2</sub> is a child of N<sub>1</sub> and a parent to T<sub>1</sub> and T<sub>2</sub>. N<sub>3</sub> is a child of N<sub>1</sub> and a parent to T<sub>3</sub> and T<sub>4</sub>. Edges are lines connecting a parent node to a child node or a parent node to a tip.



476  
477  
478  
479  
480

**Figure 2. Schematic of PhyC, Synchronous, and Continuous Tests.** Tree edges indicate:  
binary phenotype presence (pink), continuous phenotype value (rainbow),  
genotype presence (light blue), (dark blue), and (purple).

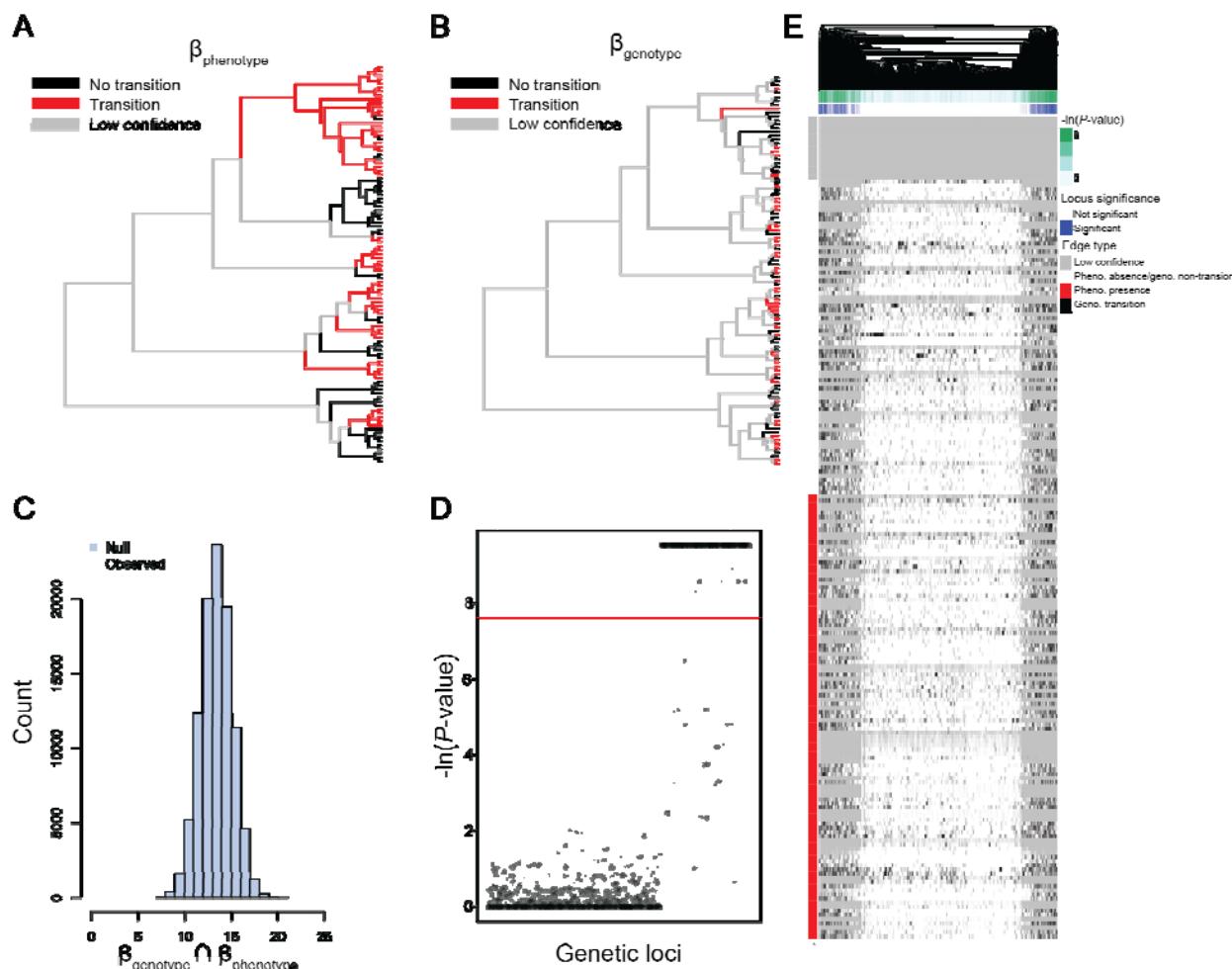


481  
482  
483  
484

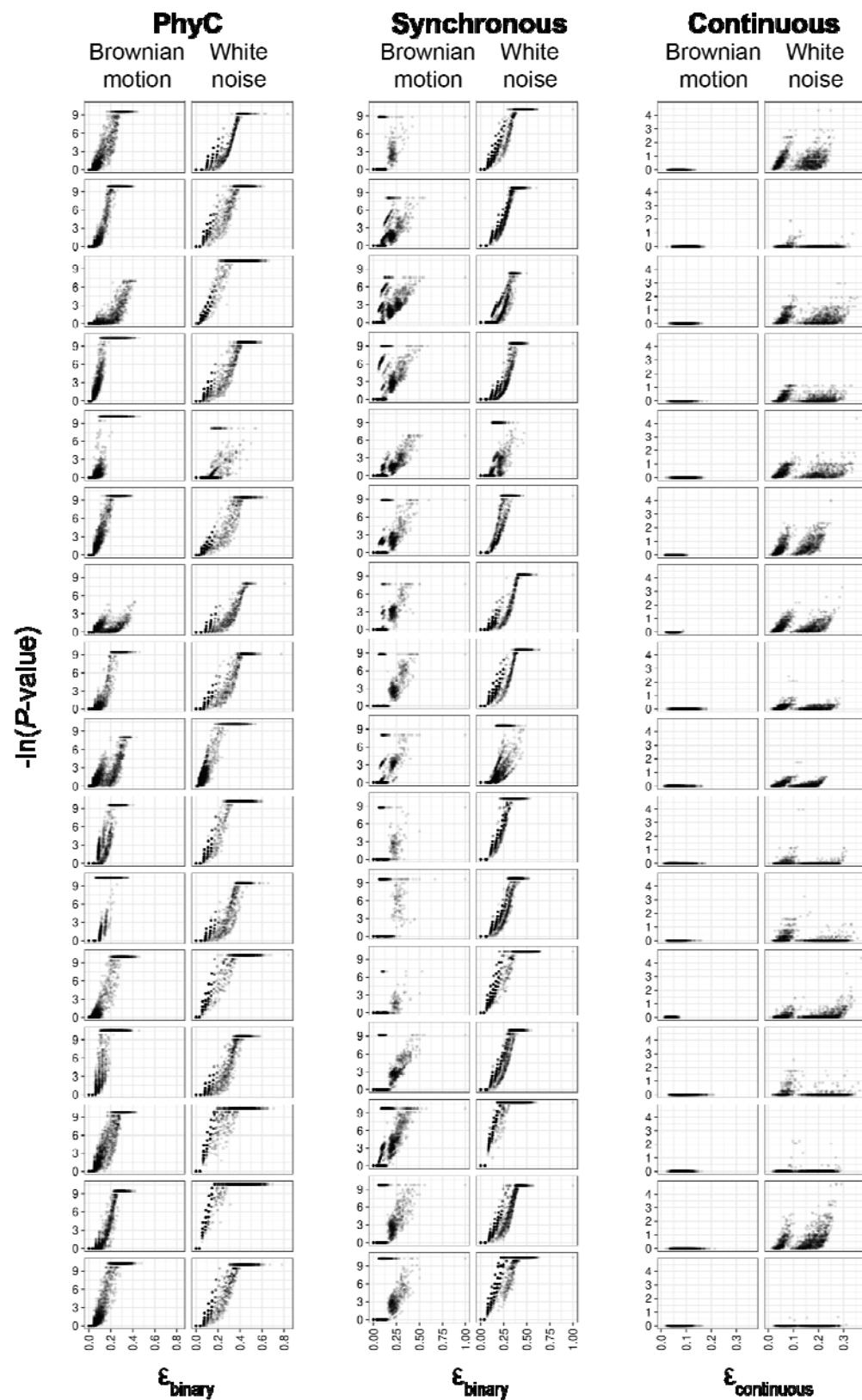
**Figure 3. Example of hogwash grouping feature.** In this case, three SNPs are found in the same gene (Gene A). No individual SNP is convergent on the tree. Hogwash performs ancestral state reconstruction on each SNP. The edges where SNP presence is inferred are colored.

485 Next, hogwash identifies the  
486 combines the three SNP  
487 When the SNPs are grouped into Gene A  
488

for each SNP (colored edges). Finally, hogwash  
together to create the Gene A (purple edges).



489  
490 **Figure 4. Example output from hogwash PhyC results from simulated data.** A) Edges with:  
491 phenotype presence ( ) in red; phenotype absent in black; low confidence in tree or  
492 low confidence phenotype ancestral state reconstruction in gray. B) Edges with: genotype  
493 mutations that arose ( ) in red; genotype mutation did not arise in black; low confidence  
494 in tree or low confidence genotype ancestral state reconstruction in gray. C) Null distribution of  
495 . D) Manhattan plot for all tested genotypes. E) Heatmap with tree edges  
496 in the rows and genotypes in the columns. The genotypes are hierarchically clustered. The  
497 genotypes are classified as being a transition edge in black or non-transition edge in white. The  
498 column annotations pertain to loci significance; green indicates the  $P$ -value while blue indicates  
499 that the  $P$ -value is more significant than the user-defined threshold. The row annotation  
500 classifies the phenotype at each edge; red indicates phenotype presence and white indicates  
501 phenotype absence. Gray indicates a low confidence tree edge; low confidence can be due to  
502 low phenotype ancestral state reconstruction likelihood, low genotype ancestral state  
503 reconstruction likelihood, low tree bootstrap value, or long edge length.  
504



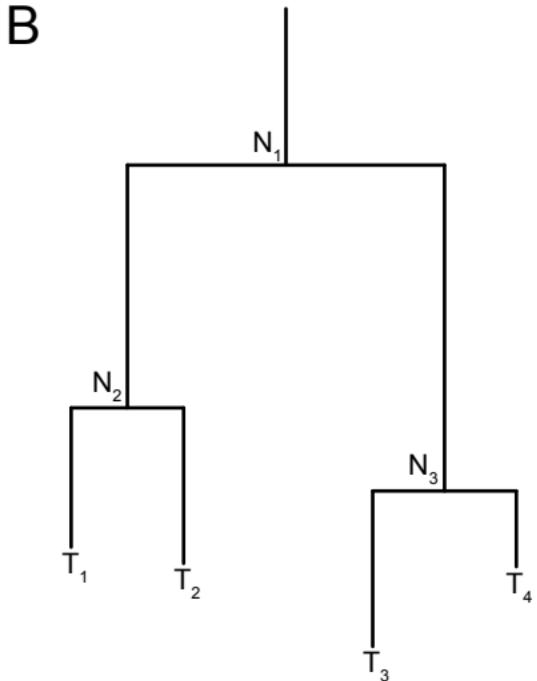
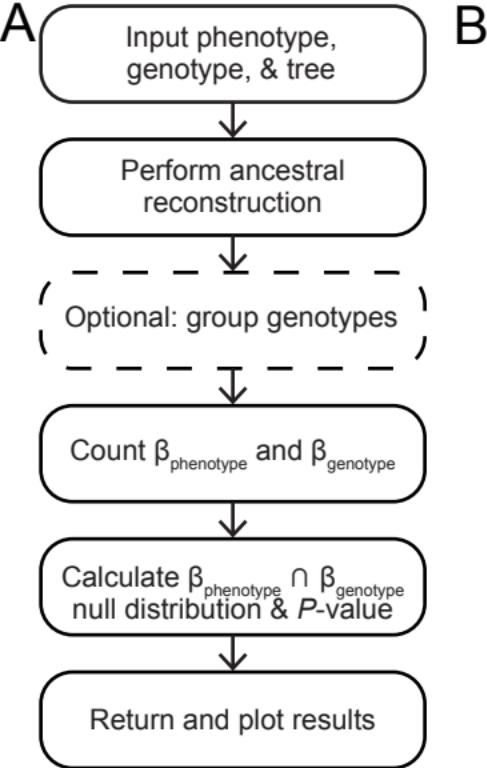
505  
506  
507  
508

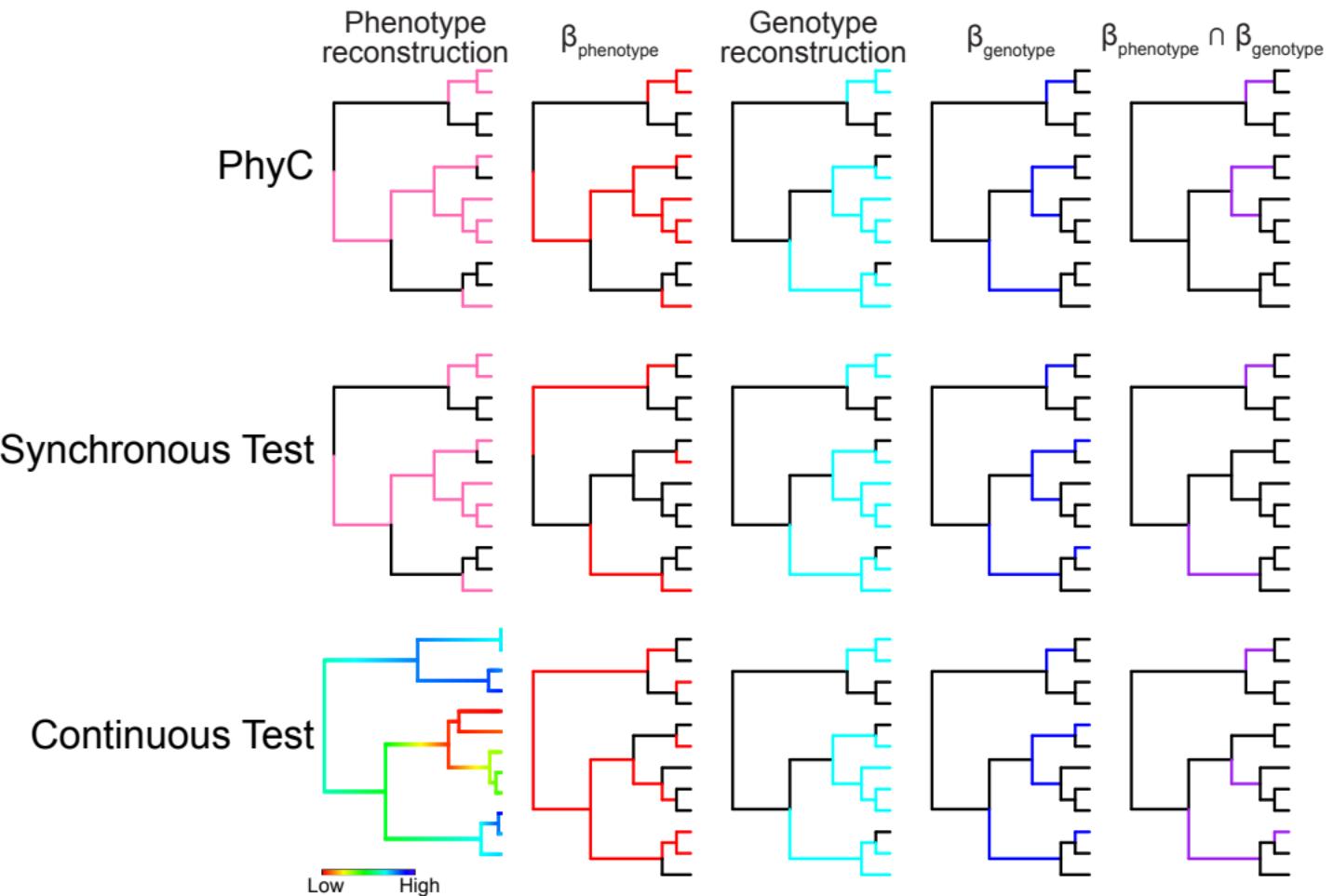
**Figure 5. High  $\epsilon$  values correlate with increased significance.** Each plot is a tree-phenotype pair. Each point represents one genotype-phenotype pair.

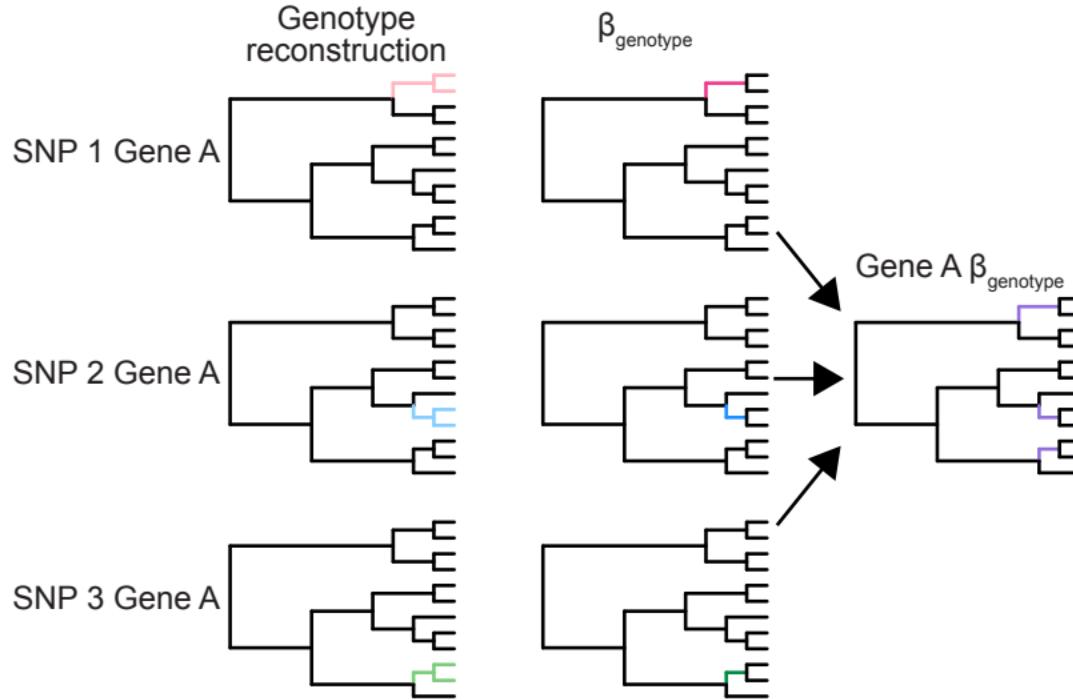
509 **TABLES**

	Phenotype	
	Brownian motion	White noise
PhyC	0.91	0.93
Synchronous Test	0.60	0.94
Continuous Test	NA	0.08

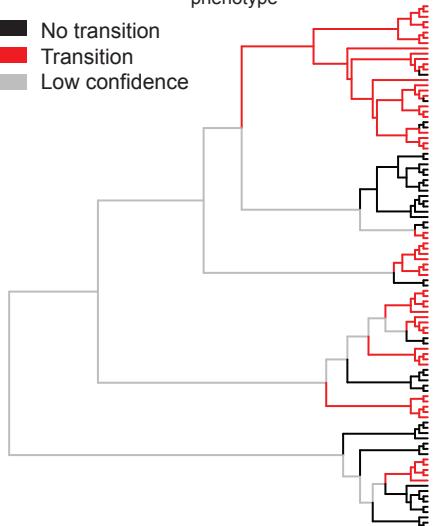
510  
511 **Table 1. Mean Spearman's rank correlation coefficient for -ln(P-value) versus from**  
512 **hogwash run on simulated data.** The  $\rho$  could not be calculated for the results from the  
513 Continuous Test on the Brownian motion phenotypes because, after multiple testing correction,  
514 all  $P$ -values are identical.  
515



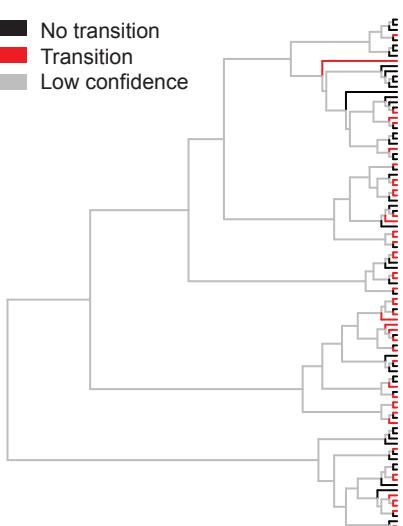




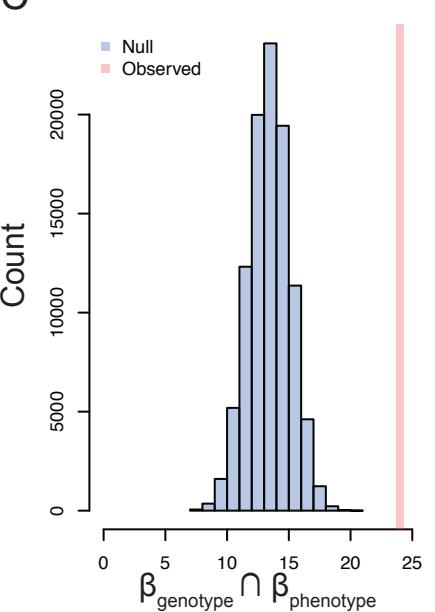
A



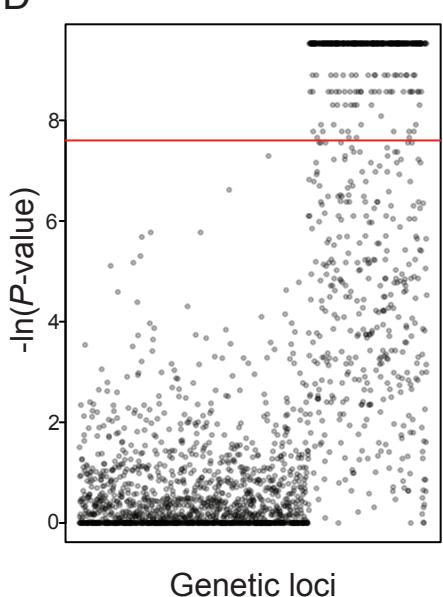
B



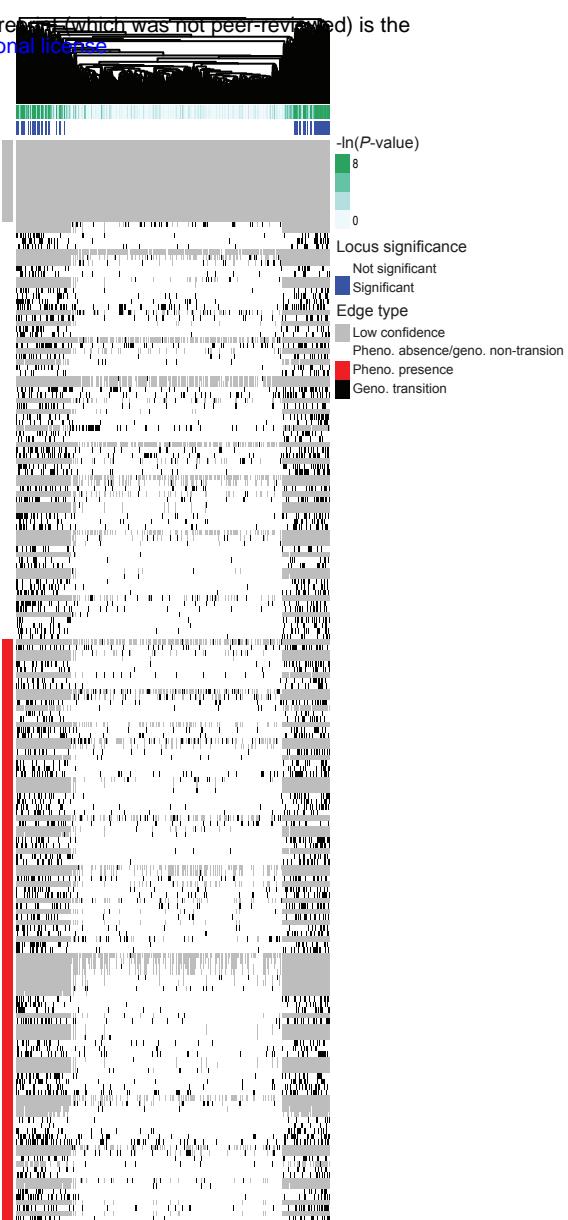
C

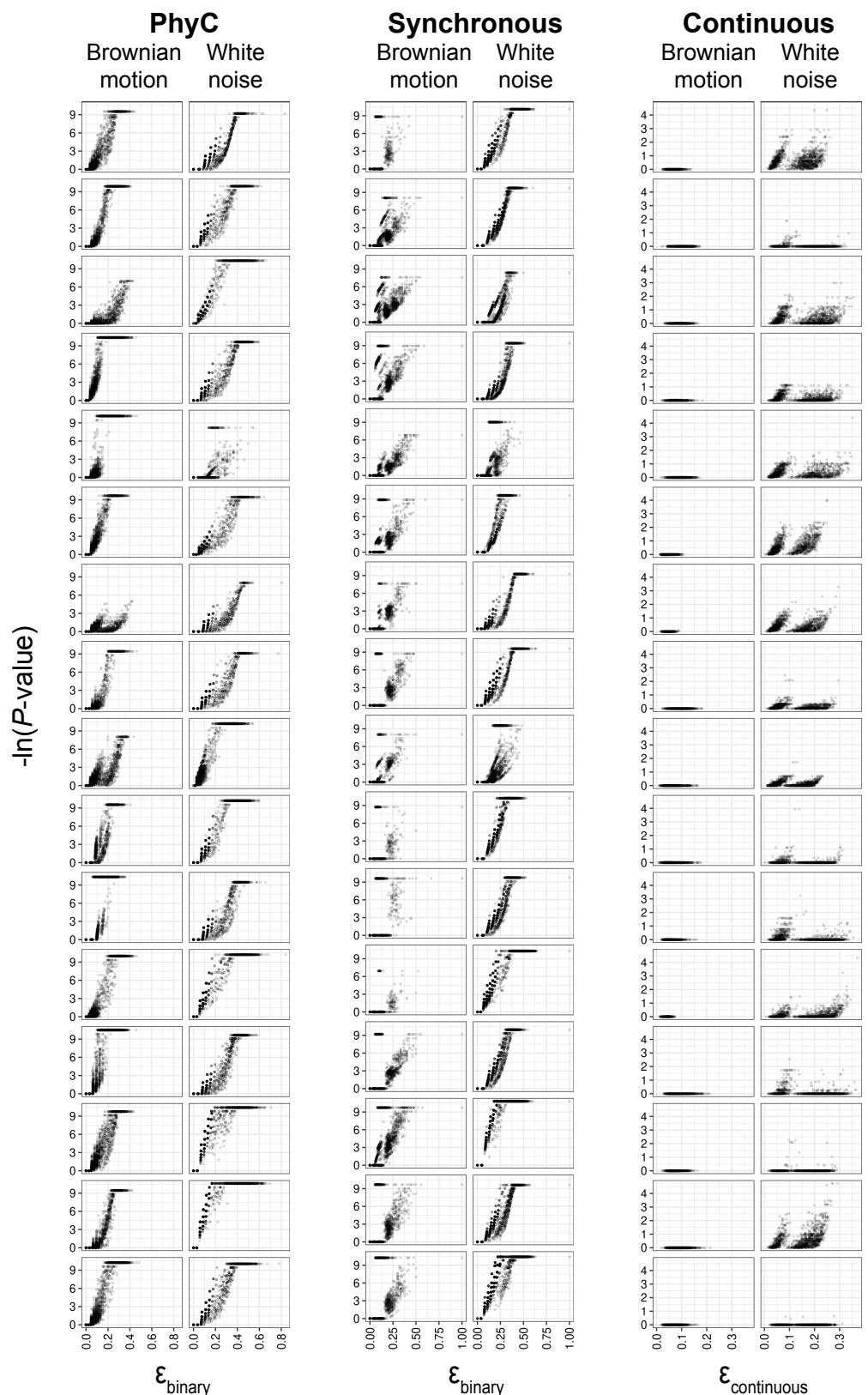


D



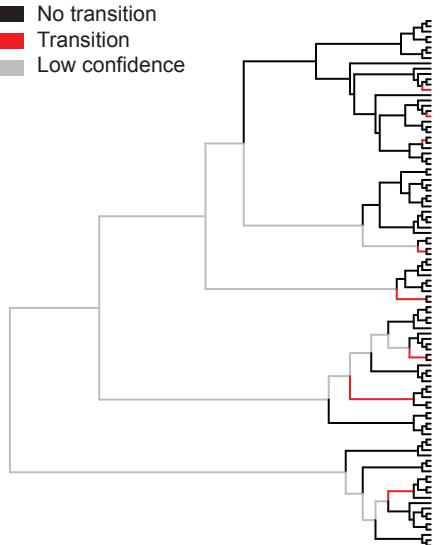
E



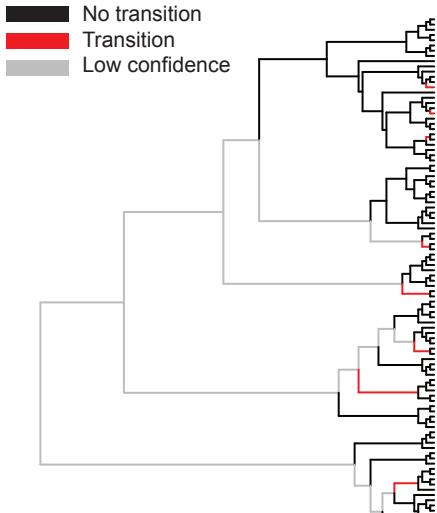


	Phenotype	
	Brownian motion	White noise
PhyC	0.91	0.93
Synchronous Test	0.60	0.94
Continuous Test	NA	0.08

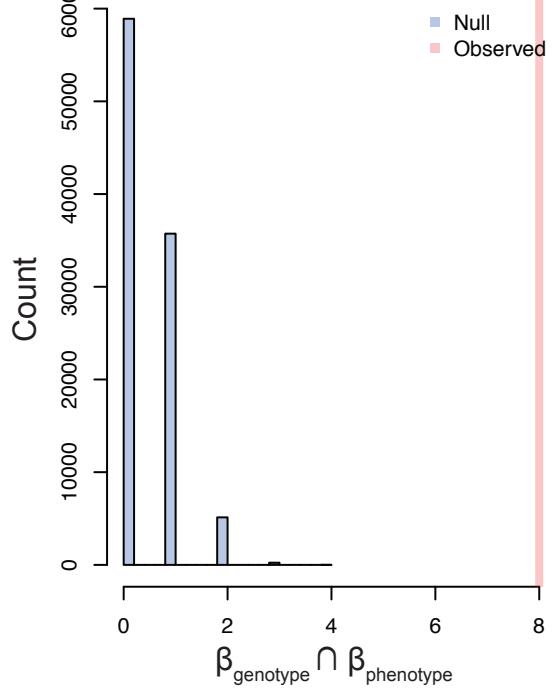
A



B



C



D

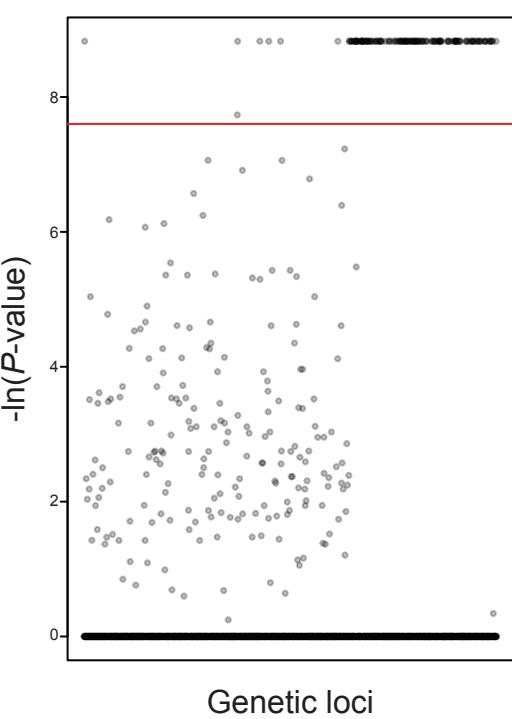
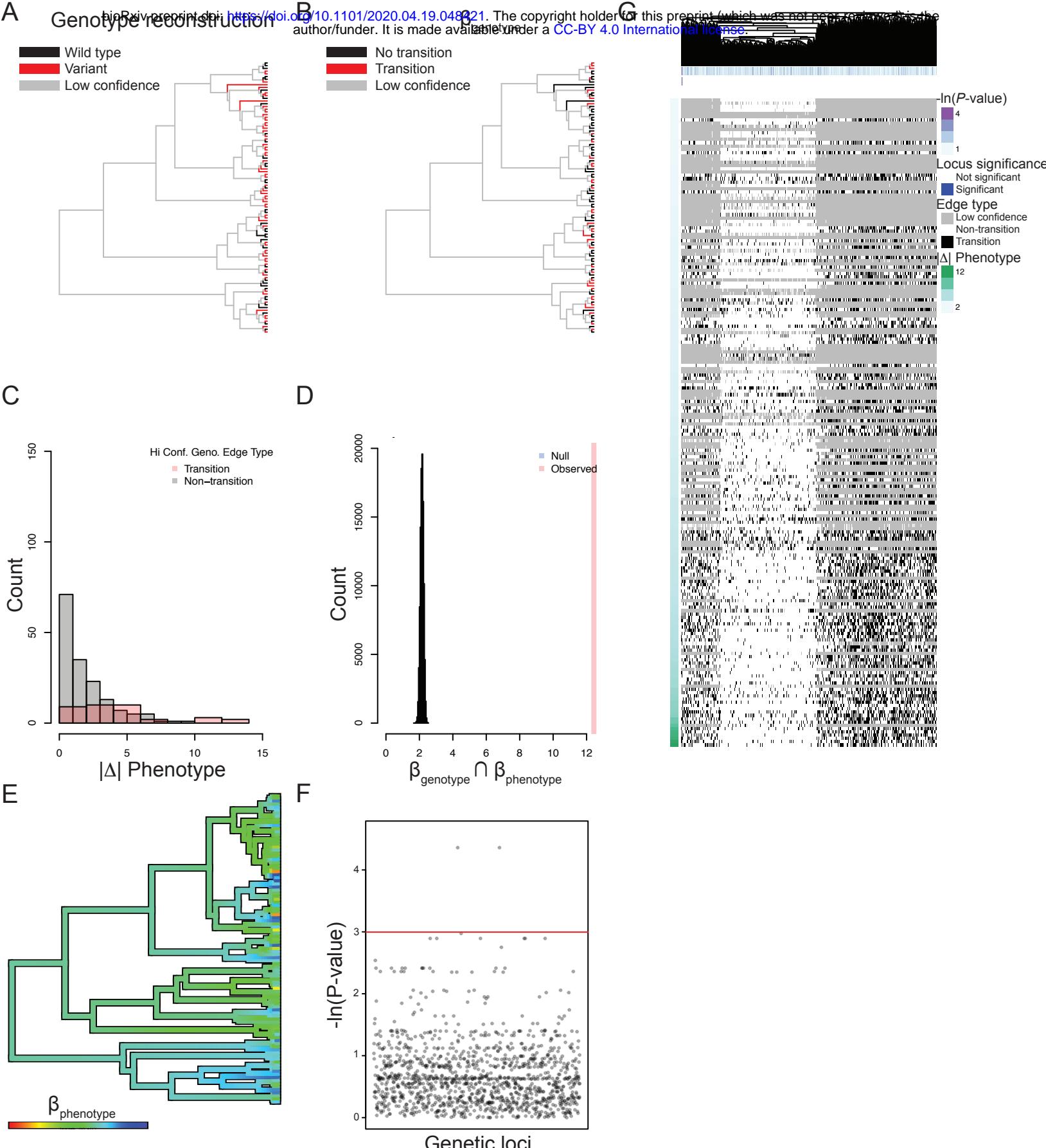


Figure S1. Example output from hogwash Synchronous Test results on simulated data

A) Ancestral reconstruction of phenotype. Phenotype transition edges ( $\beta_{\text{phenotype}}$ ) in red, phenotype non-transitions edges in black, low confidence edges in gray. B) Ancestral reconstruction of simulated genotype. Genotype transition edges ( $\beta_{\text{genotype}}$ ) in red; non-transition edges in black; low confidence edges in gray. C) Null distribution of  $\beta_{\text{genotype}} \cap \beta_{\text{phenotype}}$ ; observed value in red. D) Manhattan plot for all tested genotypes. Significance threshold indicated in red. E) Heatmap with tree edges in the rows and genotypes in the columns. The genotypes are hierarchically clustered. The genotypes are classified as being a transition edge in black or non-transition edge in white. The column annotations pertain to the  $P$ -value; green is the  $P$ -value and blue indicates that the  $P$ -value is more significant than the user-defined threshold. The row annotation classifies the phenotype edge type. Phenotype transition edges are in red and non-transition edges are in white. Gray indicates a low confidence tree edge; low confidence can be due to low genotype ancestral state reconstruction likelihood, low tree bootstrap value, or long edge length.



**Figure S2.** Example output from the Continuous Test run on simulated data.

A) Reconstruction for a simulated genotype. Wild type in black, variant presence in red, and low confidence edges in gray. B) Genotype transition edges ( $\beta_{\text{genotype}}$ ) in red; non-transition edges in black; low confidence in tree or low confidence genotype ancestral state reconstruction in gray. C) Histogram of the change in phenotype per edge for high confidence tree edges. Genotype transition edges in red and genotype non-transition edges in gray. D) Null distribution of  $\beta_{\text{genotype}} \cap \beta_{\text{phenotype}}$ ; observed value in red. E) Ancestral reconstruction of phenotype. F) Manhattan plot for all tested genotypes. The significance threshold is indicated in red. G) Heatmap with tree edges in the rows and genotypes in the columns. The genotypes are hierarchically clustered. The genotypes are classified as being a transition edge in black or non-transition edge in white. The column annotations pertain to the  $P$ -value; purple indicates the  $P$ -value and blue indicates that the  $P$ -value is more significant than the user-defined threshold. The row annotation shows the absolute value in the phenotype change per edge. Gray indicates a low confidence tree edge; low confidence can be due to low genotype ancestral state reconstruction likelihood, low tree bootstrap value, or long edge length.

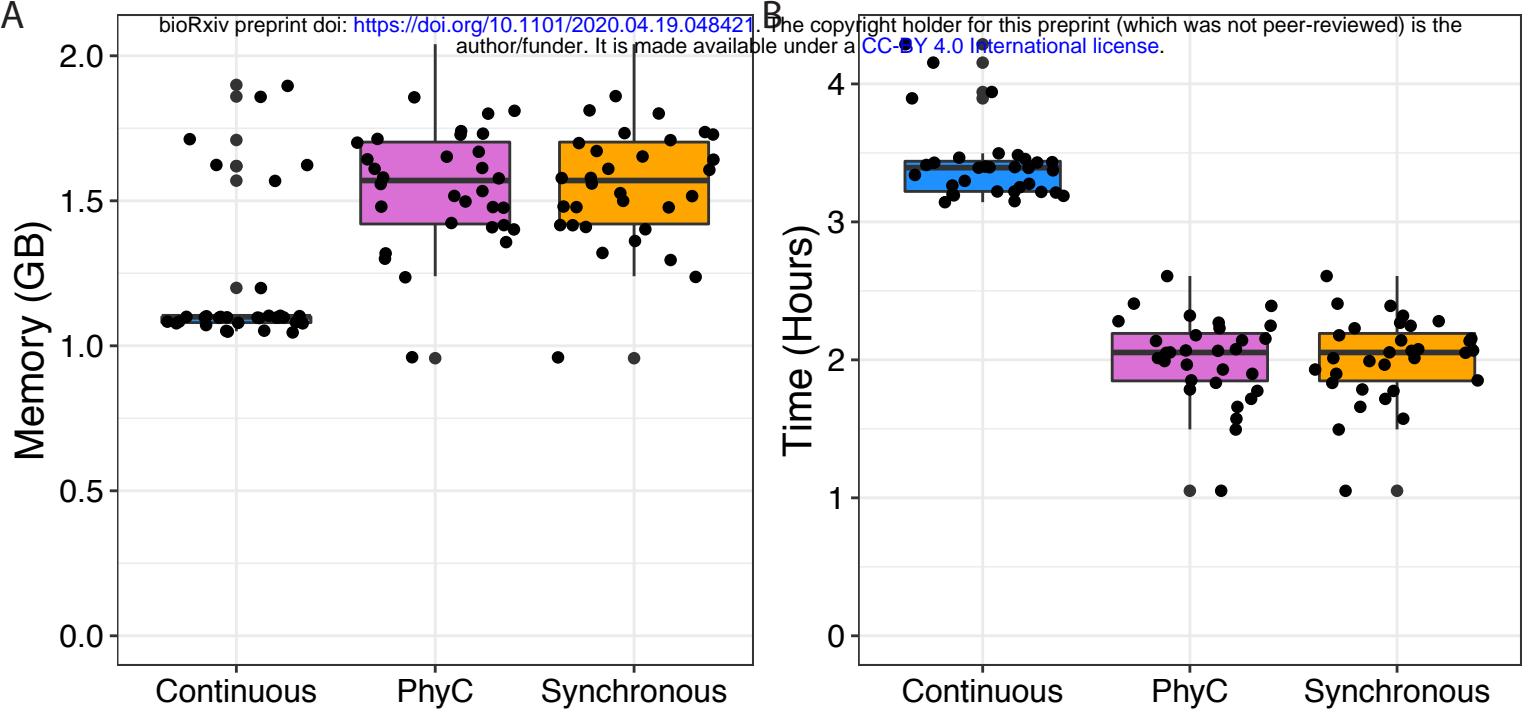


Figure S3. Memory usage and run time for hogwash on simulated data.